intel.

# Simplifying Generative AI App Development:
# Why Standards Matter
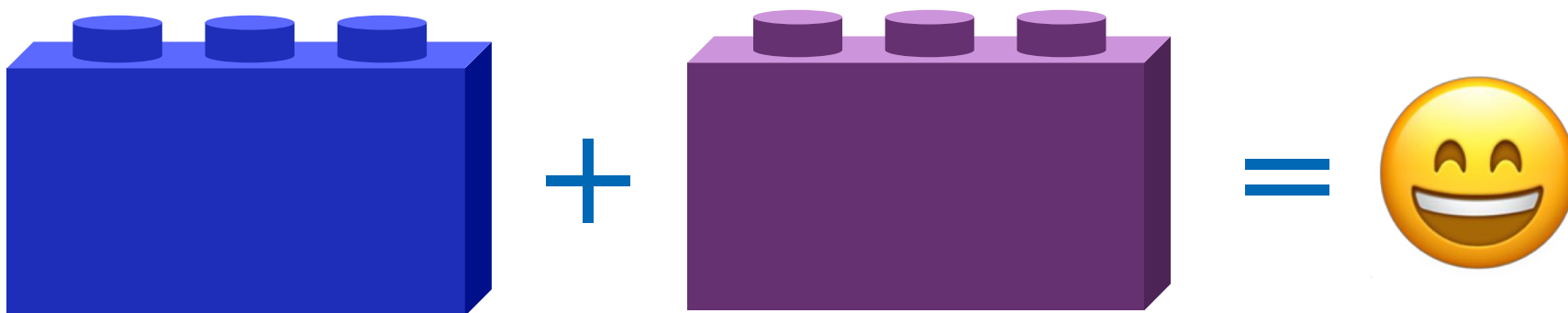
**Katherine Druckman**

Open Source Security Evangelist

**Ezequiel Lanza**
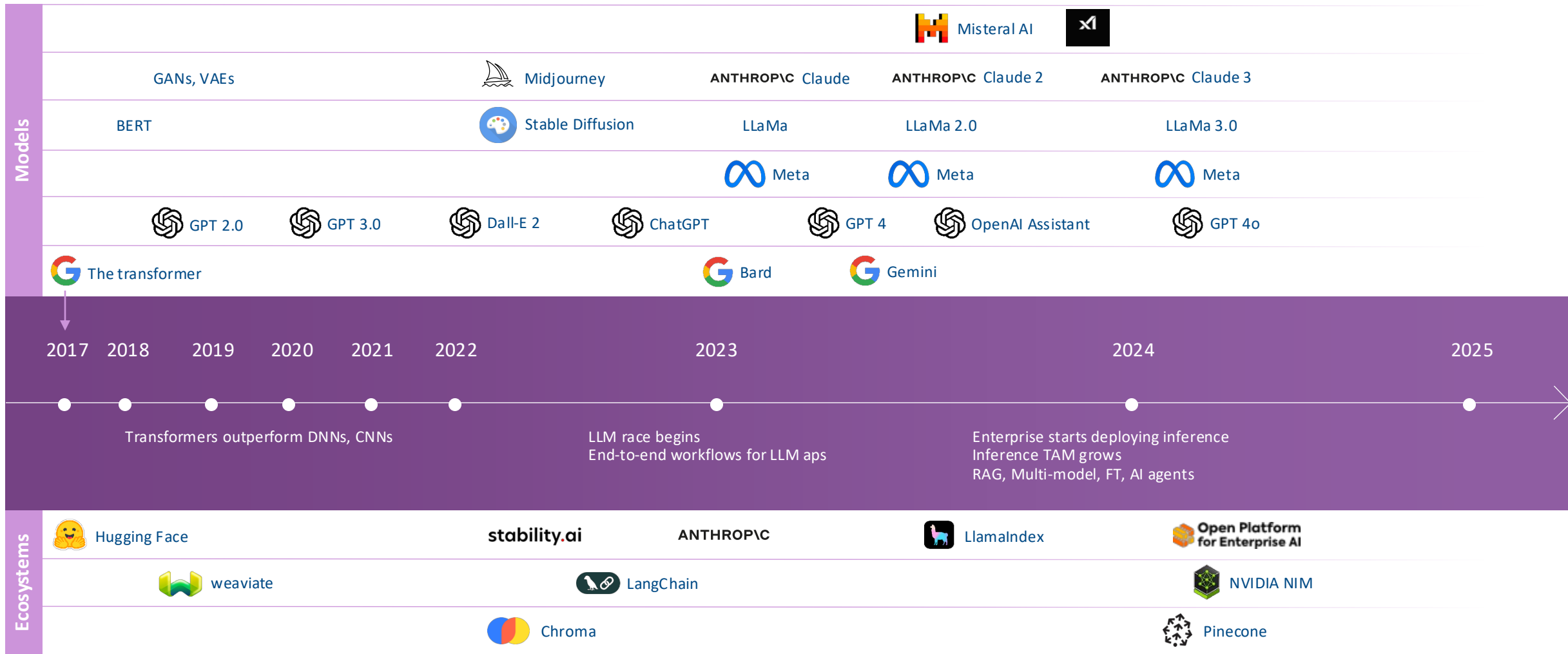
Open Source AI Evangelist

GenAI is emerging rapidly, but enterprises are struggling to realize GenAI value in production.

# Life is easier when things go together

# GenAI: transformers to AI Agents in 8 years



**Models**

| | | | | | |
|---|---|---|---|---|---|
| | | | | Misteral AI | xAI |
| GANs, VAEs | Midjourney | ANTHROP\C Claude | ANTHROP\C Claude 2 | ANTHROP\C Claude 3 | |
| BERT | Stable Diffusion | LLaMa | LLaMa 2.0 | LLaMa 3.0 | |
| | | Meta | Meta | Meta | |
| GPT 2.0 | GPT 3.0 | Dall-E 2 | ChatGPT | GPT 4 | OpenAI Assistant | GPT 4o |
| The transformer | | | Bard | Gemini | |

| 2017 | 2018 | 2019 | 2020 | 2021 | 2022 | 2023 | 2024 | 2025 |
|---|---|---|---|---|---|---|---|---|

Transformers outperform DNNs, CNNs

LLM race begins
End-to-end workflows for LLM aps

Enterprise starts deploying inference
Inference TAM grows
RAG, Multi-model, FT, AI agents

**Ecosystems**

| | | | | |
|---|---|---|---|---|
| Hugging Face | stability.ai | ANTHROP\C | LlamaIndex | Open Platform for Enterprise AI |
| weaviate | | LangChain | | NVIDIA NIM |
| | Chroma | | | Pinecone |

intel

# Massive innovation, expensive duplication, no standardization

- So many options, how do decide?

- Constantly reinventing the wheel

- Few best practices, wasted dev cycles.

- Multiple, circuitous routes to the same place
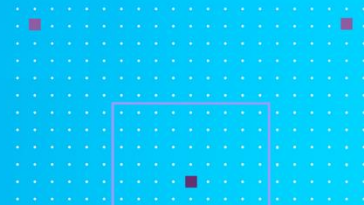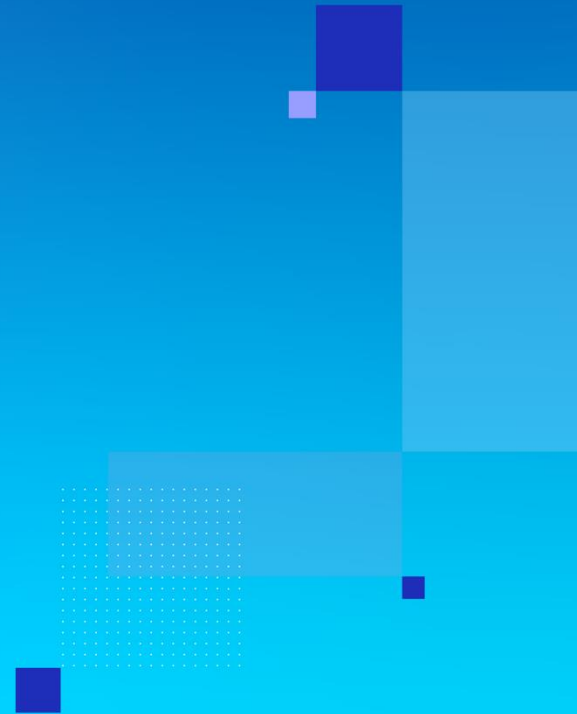
- Can we be pioneers without so much chaos?

We've created a monster!

# Without standards Gen AI RAG can be a beast!

intel

The solution

# A platform that tames the chaos

# Open source platform that organizes GEN AI chaos

THE
**LINUX**
FOUNDATION

**Open Platform for Enterprise AI**

- Composable building blocks for generative AI systems

- Integrates LLMs, data stores, and prompt engines
  RAG AI blueprint with component stack structure and
  end-to-end workflow

- Additional stacks for translation, code generation, images

- Generative AI assessment tool for performance, features,
  trustworthiness and enterprise-grade readiness

# OPEA Contributors

# What is RAG?

# LLMs don't know everything ...



But they do know a lot about their training data

# If we augment an LLM with new context …



It can generate a response to a novel topic

# RETRIEVAL

🧠

Retrieve data from a knowledge base
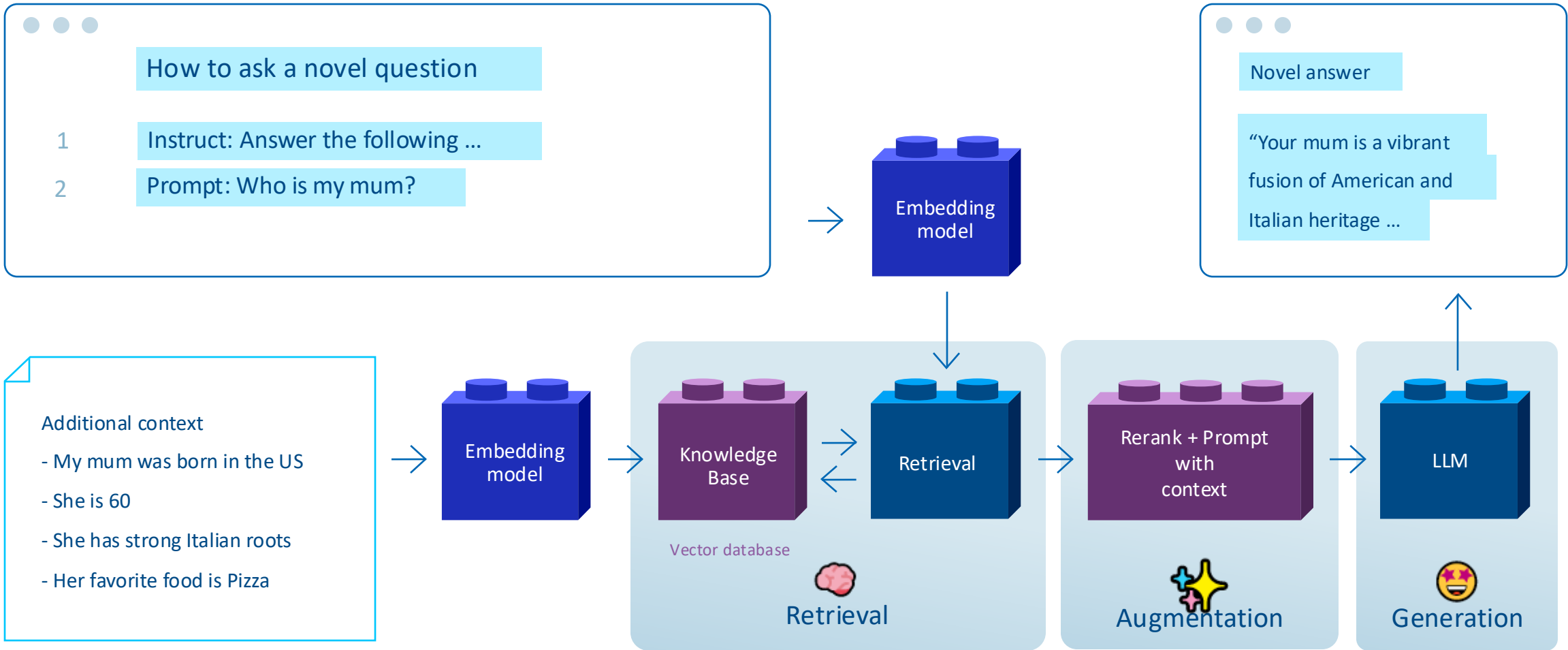
(my brain in this case!)

# AUGMENTED

✨

Give the model instructions

and context with the prompt

# GENERATION

🤩

Get a response based on the retrieved

data

# RAG process flow

# Lack of standards make Gen AI extremely complex

Ecosystem complexity for Gen AI

**Data**

- Indexing the data

- Creating embeddings (vector representations)

- Choosing storage (ISVs?)

- Preparing the data

**LLM**

- Selecting the right model

- Interacting with multiple models

- Choosing frameworks

**Deploy**

- Hardware efficiency

- Retrieval method

- Applying guardrails

- Retrieval accuratcy (reranking)

- Putting it all together
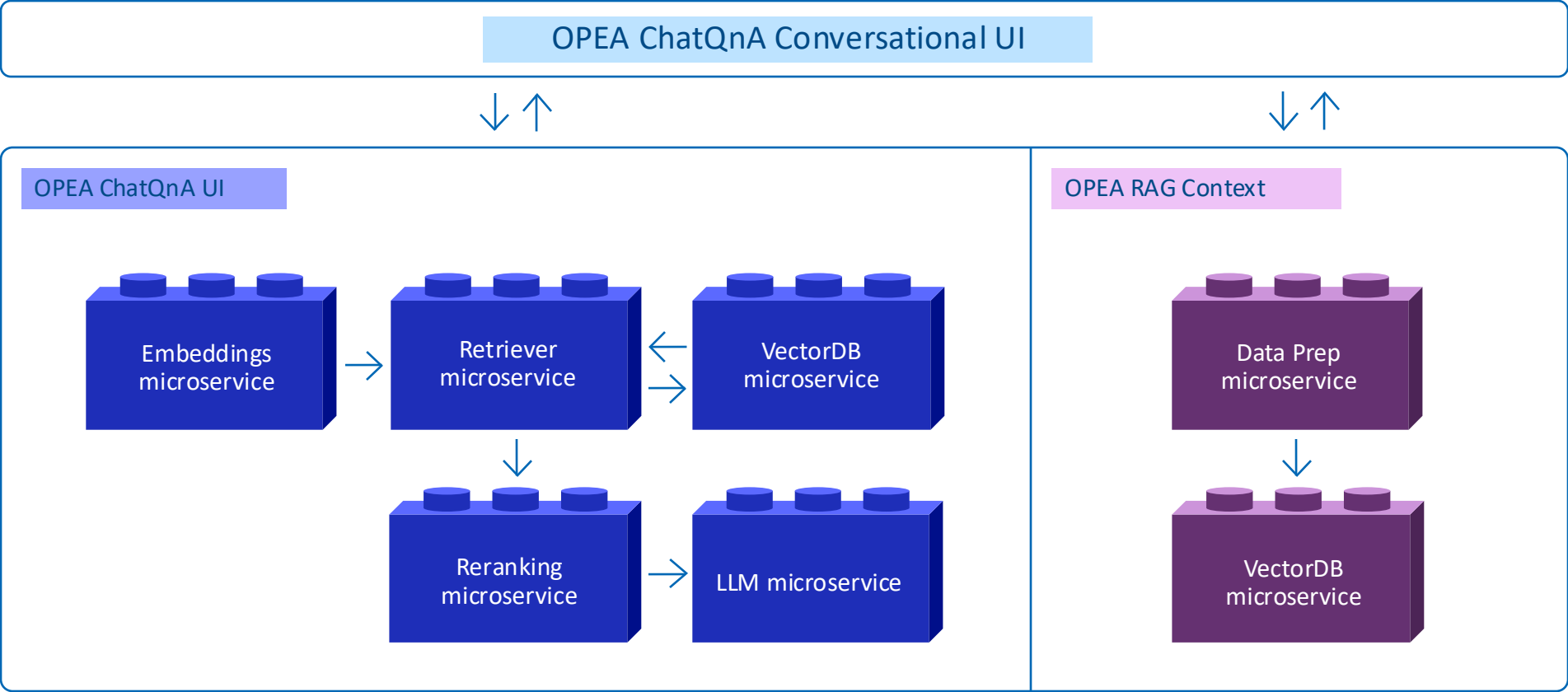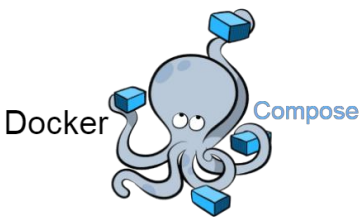
- Finding/building a megaservice

## OPEA RAG reference meets these challenges

- Composable microservices for data, LLMs, and deployment

- Plug and play with 3rd party services—point to whichever  data providers,  embedding engines, and models you choose

- Over-arching mega service in place for faster development

# DEMO

# Under the hood – OPEA ChatQnA with RAG

Docker Compose



OPEA ChatQnA Conversational UI

OPEA ChatQnA UI

OPEA RAG Context

Embeddings microservice → Retriever microservice ← VectorDB microservice

Data Prep microservice

Reranking microservice → LLM microservice

VectorDB microservice

- Chat with RAG running on microservice architecture

- Microservices can point to choice of services/APIs

Open Platform for Enterprise AI

# Embeddings microservice - Converts from text to embeddings (vectors)

opea/embedding-tei ☆0

By opea · Updated 3 days ago

IMAGE
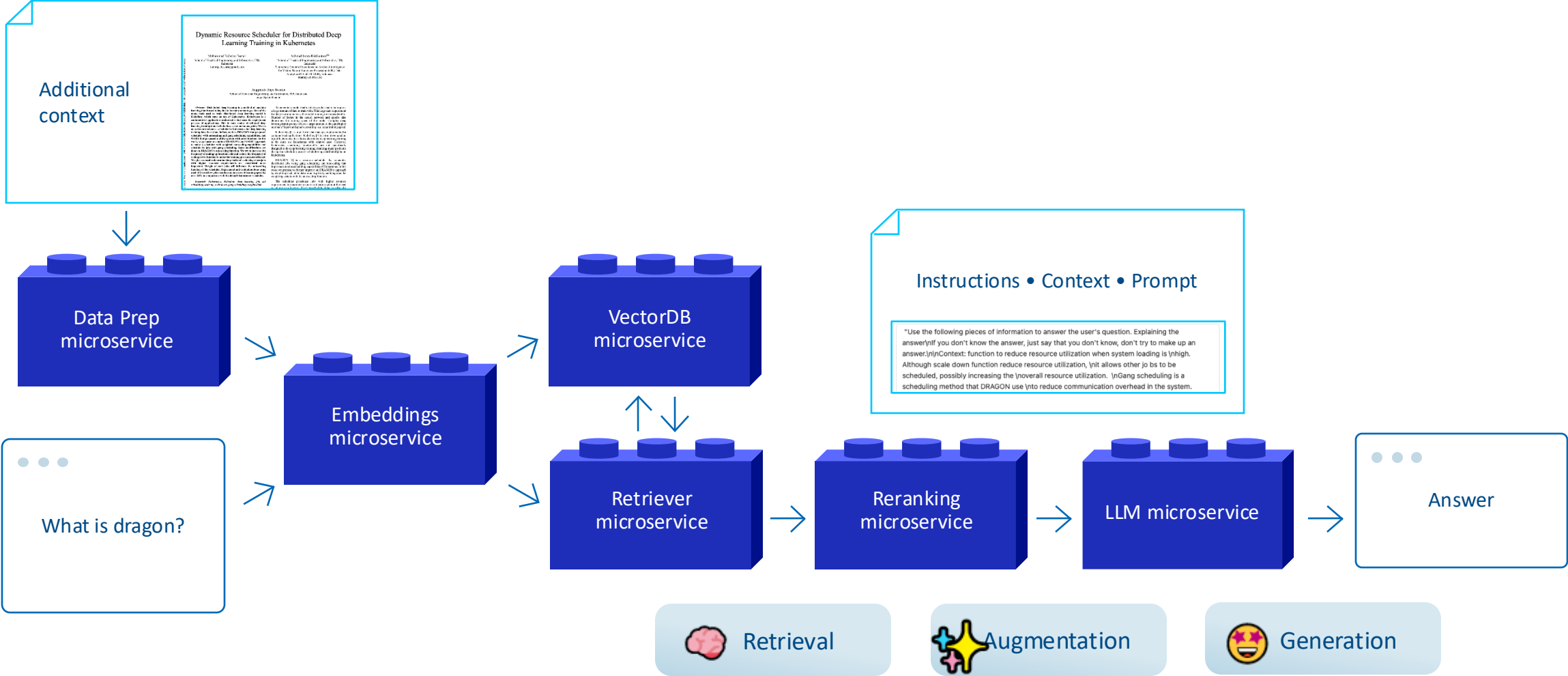
Variables

```
services:
  embedding:
    image: opea/embedding-tei:latest
    container_name: embedding-tei-server
    ports:
      - "6000:6000"
    ipc: host
    environment:
      http_proxy: ${http_proxy}
      https_proxy: ${https_proxy}
      TEI_EMBEDDING_ENDPOINT: ${TEI_EMBEDDING_ENDPOINT}
      LANGCHAIN_API_KEY: ${LANGCHAIN_API_KEY}
    restart: unless-stopped
```
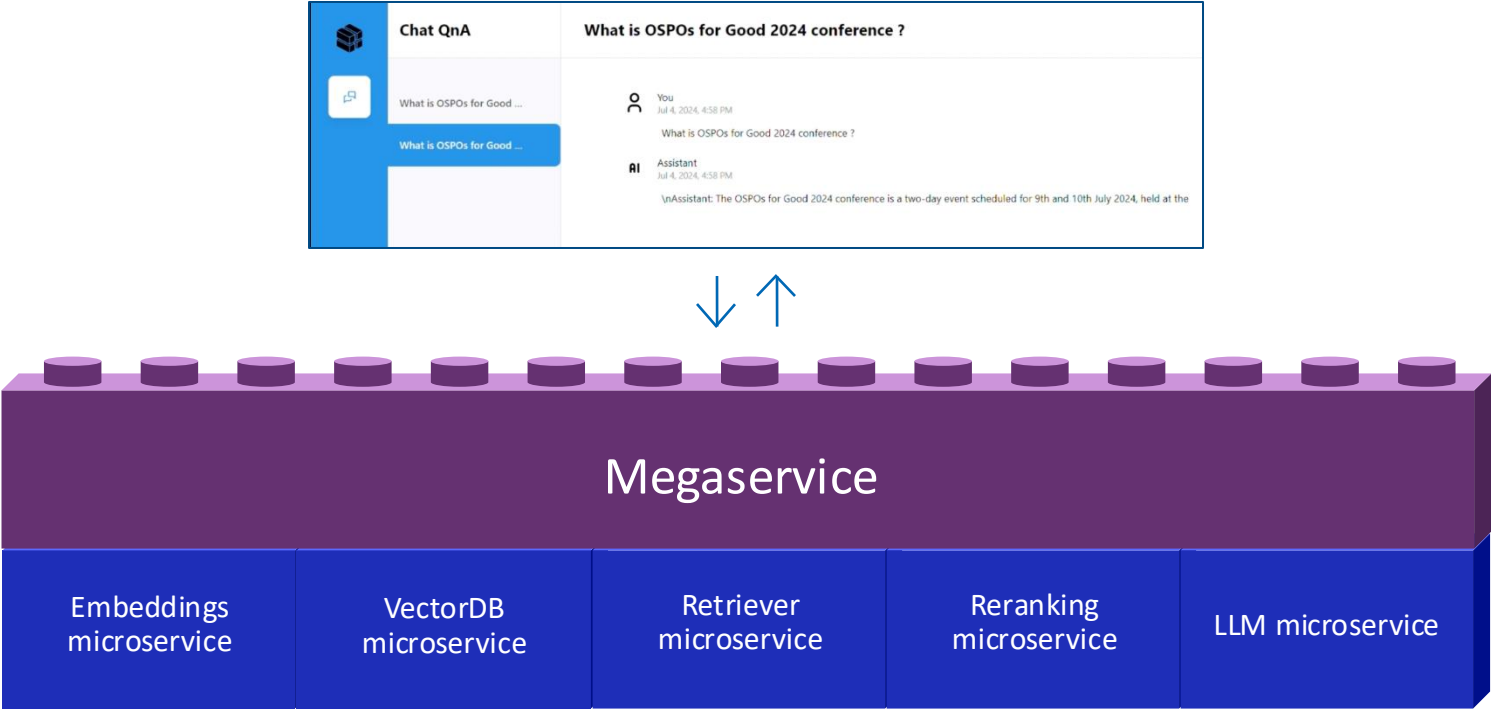
# Vector for "Who is my mum?"

```
curl localhost:$your_port/embed    -X POST    -d
'{"inputs":"Who is my mum?"}'    -H 'Content-Type:
application/json'
```

```
[[0.018129222,0.0030305043,-0.049874727,-
0.035030127,0.014229514,-0.023594731………..-
0.033771276,-0.0009737879,-0.006777766,-
0.058805678,0.011158894,0.012094927,0.01805739,
0.054448325,-0.032204594,0.049175356]]
```

# How OPEA ChatQnA answers on context



Additional context

Dynamic Resource Scheduler for Distributed Deep Learning Training in Kubernetes

Data Prep microservice

Embeddings microservice

VectorDB microservice

What is dragon?

Retriever microservice

Reranking microservice

LLM microservice

Answer

**Instructions • Context • Prompt**

"Use the following pieces of information to answer the user's question. Explaining the answer\nIf you don't know the answer, just say that you don't know, don't try to make up an answer.\n\nContext: function to reduce resource utilization when system loading is \nhigh. Although scale down function reduce resource utilization, \nit allows other jo bs to be scheduled, possibly increasing the \noverall resource utilization. \nGang scheduling is a scheduling method that DRAGON use \nto reduce communication overhead in the system.

🧠 Retrieval    ✨ Augmentation    🤩 Generation

# OPEA ChatQnA Megaservice

# CODE

intel.

# Check out ChatQnA on GitHub

# RAG conclusions and challenges

The good

- RAG is a great approach to having context-based answers

- OPEA makes the building blocks easy to deploy

The challenges

- Dealing with a large amount of data (Play with different similar search techniques)

- Different types of data (Images/Tables ) – Intel Tech on Medium articles

- Is it all about indexing? Explore Retrieval augmented text-to-SQL generation

THE LINUX FOUNDATION

Open Platform for Enterprise AI

Join us!

There is more—a lot more—you can build!

# OPEA* Roadmap

E2E Examples and Ready-to-deploy Reference Flows

Open Models for OPEA Integration

Optimized Compilers and Toolchains for OPEA

Functional Components with Tools or Microservice

Retrieval Augmented Generation (RAG) – Standardized modular, heterogenous RAG pipelines for enterprise

Finetuning Pipeline

Multimodality (Image/Video)

AI Agent (Single and Multi)

Evaluation and Benchmarking

Vertical Use Case Working Groups, focused on industry challenges

*Preliminary

# Help build a better OPEA, contribute!

# Participate in OPEA



Visit OPEA.dev

- Join a Working Group
- Bring Enterprise AI use cases
- Contribute code, docs, projects, blueprints, and more
- Provide feedback
- Evangelism and Promotion of OPEA in YOUR communities, events

# Scan for tools, sources, and resources

Visit the Intel Open Ecosystem Community and

Evangelism GitHub page

- A PDF of this presentation

- Guides and community resources

- Links to articles and source material

## Where to find Katherine

- LinkedIn: /katherinedruckman
- Fediverse: @katherined@reality2.social
- Twitter/X: @katherined

## Where to find Ezequiel

- LinkedIn: /ezelanza
- Fediverse: @eze_Lanza@mstdn.ca
- Twitter/X: @eze_lanza

Open Platform for Enterprise AI

# Thank you!